

# Banking retail consumer finance data generator - credit scoring data repository

Karol Przanowski

email: [kprzan@interia.pl](mailto:kprzan@interia.pl)

url: <http://kprzan.w.interia.pl>

## Abstract

This paper presents two cases of random banking data generators based on migration matrices and scoring rules. The banking data generator is a new hope in researches of finding the proving method of comparisons of various credit scoring techniques. There is analyzed the influence of one cyclic macro-economic variable on stability in the time account and client characteristics. Data are very useful for various analyses to understand in the better way the complexity of the banking processes and also for students and their researches. There are presented very interesting conclusions for crisis behavior, namely that if a crisis is impacted by many factors, both customer characteristics: application and behavioral; then there is very difficult to indicate these factors in the typical scoring analysis and the crisis is everywhere, in every kind of risk reports.

**Key words:** credit scoring, crisis, banking data generator, retail portfolio.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Detailed description of data generator</b>	<b>6</b>
2.1	The main options . . . . .	6
2.2	Production dataset . . . . .	6
2.3	Transaction dataset . . . . .	7
2.4	Inserting the Production dataset into the Transaction dataset	7
2.5	Analytical Base Table – ABT dataset . . . . .	8
2.6	Migration matrix adjustment . . . . .	9
2.7	Iteration step . . . . .	9
2.8	Default definition . . . . .	10
2.9	Portfolio segmentation and risk measures . . . . .	11
<b>3</b>	<b>General theory</b>	<b>12</b>
3.1	The main assumption and definition . . . . .	12
3.2	Open questions . . . . .	12
<b>4</b>	<b>Two case studies</b>	<b>13</b>
4.1	Common parameters . . . . .	13
4.2	The first case study – <b>unstable application characteristic</b> – APP . . . . .	14
4.3	The second case study – <b>unstable behavioral characteris-</b> <b>tic</b> – BEH . . . . .	16
4.4	Stability problem . . . . .	16
4.5	Various types of risk measures . . . . .	18
4.6	Implementation . . . . .	18
<b>5</b>	<b>Conclusions</b>	<b>18</b>

# 1 Introduction

Currently predictive models and especially credit scoring models are very popular in management of banking processes [1]. It is a typical that risk scorecards are always used in credit acceptance process to optimize and control the risk. Various forms of behavioral scorecards are also used for management of repeat business and also for PD models in Basel RWA (Risk Weighted Assets) calculation [2]. It is a kind of phenomenon that a list of about 10 account or client characteristics can predict their future behavior, their style of payments and their delinquency.

One can say the trivial fact scorecards are useful and methodology is well known, but on the other hand still credit scoring can be developed and new techniques should be tested. The main problem today is that there is not defined the general testing idea of new methods and techniques, there is no proving method of their correctness. Many good articles are prepared based on one particular case study, on one example of real data coming from one or a few banks [3], [4] and [5]. From a theoretical point of view, even there are presented good results and very correct arguments to suggest choosing one method than another, it is the prove only on that particular data is indicated the difference, but nobody can prove it for other data, nobody can guarantee the correctness for all cases.

There are also other important reasons why real banking data are not available globally and cannot be used by everyone analysts, like legal constraints or too fresh new products with too short data history. These two factors suggest finding a quite another approach for predictive modeling testing in banking usage.

It is a very good idea to start developing two parallel ways: real data and random-simulated data approaches. The second one even cannot replace real data it can be very useful to understand in the better way relations among various factors in data, to imagine a complexity of the process and can be a trial to create more general class of semi-real data.

Let be considered some advantages of randomly generated data:

1. Today many analysts try to understand and to analyze the last crisis [6], among other things they develop methods of indicating risk stable in the time sub-portfolios. Topic is not easy and cannot be solved by typical predictive models based on target variable like in the case of default

risk. The notion of stability cannot be defined for every particular account or client, one cannot say that account is stable, only the set of accounts can be tested, so that technique should be developed by quite different method than typical predictive modeling with target variable. It can be formulated by a simple conclusion: the more accounts the more robust stability testing. In the random data generator can be tested various scenario to see and to better understand the problem.

2. Scoring Challenges or Scoring Olympic Games. From time to time there are organized by different environments contests to find good modelers or to test new techniques. Sometimes data are taken from too real case. Too real means, that some real processes are not predictable, because they are influenced by many immeasurable factors. Even if scoring models are used in practice also in these cases it is not a good idea to use that data for contest. The best solution and the best fairly is to use random data generator process directly predictable.
3. Reject inference area [4]. Still that topic needs development. Random data can be generated also for, in the reality, rejected cases for testing, so it can be used for better estimation of risk on blank areas and better experience.
4. Today there are two or more techniques of scorecard building [7]. It needs to make some comparisons, to make some analysis to define recommendations: where and what conditions suggest to use one than another method. The same case can be applied for different variable selection methods.
5. Product profitability, bad debts and cut-offs. On random data all mentioned notions can be tested and analysts experience can be broadened.
6. Random data can also be very important factor in the topic of data standardization or the idea of auditing. Let imagine that there are prepared all ready run software tools for MIS (Management Information Systems) and KPI (Key Performance Indicators) reporting on the generic data structure firstly uploaded by random data. Then auditing of all another data will be minimized by only the upload data process.

Simulation data are used in many areas, for example it is very useful in research of telecommunication network by the system like OPNET [8]. Also

there are developed simulated data in the banking area by [9] and [10].

The simplest retail consumer finance portfolio is the fixed installment loan portfolio. Here process can be simplified by the following assumptions:

- for all accounts one due date in the middle of the month is defined (every 15th),
- every client has only one credit,
- client can pay whole one installment, a few installments or pay nothing, two events only: payment or missing payment,
- there are measured delinquency on state: end of month by indicated the number of due installments,
- all customer and account properties are randomly generated by defined proper random distributions,
- if the number of due installments attain 7 (180 past due days) the process is stopped and account is marked by bad account status, next collection steps are omitted,
- if number of paid installments attains the number of all installments then the process is stopped and account is marked by closed account status,
- payments or missing payments are determined by three factors: score calculated on account characteristics, migration matrix and adjustment of that matrix by one cycle time macroeconomic variable,
- score is calculated for every due installments group separately. In more general case there can be defined different score for every status: due installments 0, 1, ..., and 6.

It is a good circumstance to emphasize that risk management today has very good tools for risk control, even if the crisis has come and was not predicted in the correct way, it could be indicated very quickly. It seems that the best of risk control tools is the migration matrix reporting.

The goal of that paper can be also formulated in the following way: to create random data with the condition to obtain the same results like observed in the reality by typical reporting like migration matrix, flow-rates or roll-rates and vintage or default rates.

## 2 Detailed description of data generator

### 2.1 The main options

All data are generated from starting date  $T_s$  to ending  $T_e$ .

The migration matrix  $M_{ij}$  (transition matrix) is defined as a percent of transition after one month from due installments  $i$  to due installments  $j$ .

There is one macro-economic variable dependent only on a time by the formula:  $E(m)$ , where  $m$  is a number of month from  $T_s$ . It should satisfy the simple condition:  $0.01 < E(m) < 0.9$ , because it is used as an adjustment of migration matrix, so it influences on the risk; in some months produces slightly greater one and in some months lower.

### 2.2 Production dataset

The first dataset contains all applications with all available customer characteristics and credit properties.

Customer characteristics (application data):

- Birthday -  $T_{Birth}$  - with the distribution  $D_{Birth}$
- Income -  $x_{Income}^a$  -  $D_{Income}$
- Spending -  $x_{Spending}^a$  -  $D_{Spending}$
- Four nominal characteristics -  $x_{Nom_1}^a, \dots, x_{Nom_4}^a$  -  $D_{Nom_1}, D_{Nom_2}, \dots, D_{Nom_4}$ , in practice they can represent variables like: job category, marital status, home status, education level, or others.
- Four interval characteristics -  $x_{Int_1}^a, \dots, x_{Int_4}^a$  -  $D_{Int_1}, D_{Int_2}, \dots, D_{Int_4}$ , represent variables like: job seniority, personal account seniority, number of households, housing spending or others.

Credit properties (loan data):

- Installment amount -  $x_{Inst}^l$  - with the distribution  $D_{Inst}$
- Number of installments -  $x_{N_{inst}}^l$  -  $D_{N_{inst}}$

- Loan amount –  $x_{Amount}^l = x_{Inst}^l \cdot x_{N_{inst}}^l$
- Date of application (year, month) –  $T_{app}$
- Id of application

The number of rows per month is generated based on the distribution  $D_{Applications}$ .

## 2.3 Transaction dataset

Every row contains the following information (transaction data):

- Id of application
- Date of application (year, month) –  $T_{app}$
- Current month –  $T_{cur}$
- Number of due installments (number of missing payments) –  $x_{n_{due}}^t$
- Number of paid installments –  $x_{n_{paid}}^t$
- Status –  $x_{status}^t$  - Active (A) - is still not paid, Closed (C) is paid, or Bad (B) – when  $x_{n_{due}}^t = 7$
- Pay days –  $x_{days}^t$  – number of days from the interval  $[-15, 15]$  before or after due date in a current month when payment was done, if there is missing payment, then pay days are also missing.

## 2.4 Inserting the Production dataset into the Transaction dataset

Every month of the Production dataset updates the Transaction dataset with the following formulas:

$$T_{cur} = T_{app}, \quad x_{n_{due}}^t = 0, \quad x_{n_{paid}}^t = 0, \quad x_{status}^t = A, \quad x_{days}^t = 0.$$

It is the process of inserting starting points of new accounts.

## 2.5 Analytical Base Table – ABT dataset

History of payments for every account is dependent on behavioral data, on behavior of previous payments. It is, of course, the assumption of that data generator.

There are many ideas of behavioral characteristics creation. There are presented the simple methods to consider the last available states and to indicate their evaluations in the time. All data are prepared in ABT datasets, the notion Analytical Base Table is used by SAS Credit Scoring Solution [11].

Let set current date  $T_{cur}$  as a fixed value. Actual states are calculated for that date by the formulas (actual data):

$$\begin{aligned}
x_{days}^{act} &= x_{days}^t + 15, \\
x_{n_{paid}}^{act} &= x_{n_{paid}}^t, \\
x_{n_{due}}^{act} &= x_{n_{due}}^t, \\
x_{utl}^{act} &= x_{n_{paid}}^t / x_{N_{inst}}^l, \\
x_{dueutl}^{act} &= x_{n_{due}}^t / x_{N_{inst}}^l, \\
x_{age}^{act} &= years(T_{Birth}, T_{cur}), \\
x_{capacity}^{act} &= (x_{Inst}^l + x_{Spending}^a) / x_{Income}^a, \\
x_{dueinc}^{act} &= (x_{n_{due}}^t \cdot x_{Inst}^l) / x_{Income}^a, \\
x_{loaninc}^{act} &= x_{Amount}^l / x_{Income}^a, \\
x_{seniority}^{act} &= T_{cur} - T_{app} + 1,
\end{aligned}$$

where  $years()$  calculates the difference between two dates in years.

Let consider two time series of pay days and due installments for the last 11 months from fixed current date by the formulas:

$$\begin{aligned}
x_{days}^{act}(m) &= x_{days}^{act}(T_{cur} - m), \\
x_{n_{due}}^{act}(m) &= x_{n_{due}}^{act}(T_{cur} - m),
\end{aligned}$$

where  $m = 0, 1, \dots, 11$ .

The characteristics indicated the evaluation in the time can be calculated by the formulas:

If every elements of time series for the last  $t$ -months are available then (behavioral data):

$$\begin{aligned}
x_{days}^{beh}(t) &= (\sum_{m=0}^{t-1} x_{days}^{act}(m)) / t, \\
x_{n_{due}}^{beh}(t) &= (\sum_{m=0}^{t-1} x_{n_{due}}^{act}(m)) / t,
\end{aligned}$$



where  $t = 3, 6, 9, 12$ .

If not all elements of time series are available then (missing imputation formulas):

$$\begin{aligned} x_{days}^{beh}(t) &= 15, \\ x_{ndue}^{beh}(t) &= 2. \end{aligned} \quad (2.1)$$

In other words behavioral variables represent average states for last 3, 6, 9 or 12 months. Without any problem user can add many other variables by replacing average statistic by another like MAX, MIN or other.

## 2.6 Migration matrix adjustment

Macro-economic variable  $E(m)$  influenses on the migration matrix by the formula:

$$M_{ij}^{adj} = \begin{cases} M_{ij}(1 - E(m)) & \text{for } j \leq i, \\ M_{ij} & \text{for } j > i + 1, \\ M_{ij} + \sum_{k=0}^i E(m)M_{ik} & \text{for } j = i + 1. \end{cases}$$

## 2.7 Iteration step

That step is running to generate next month of transactions, from  $T_{cur}$  to  $T_{cur} + 1$ . In every month some accounts are new, then the Transaction dataset is only updated by the ideas described in the subsection 2.4. Some accounts change the status by the formula:

$$x_{status}^t = \begin{cases} C & \text{when } x_{n_{paid}}^{act} = x_{N_{inst}}^l, \\ B & \text{when } x_{ndue}^{act} = 7, \end{cases}$$

and these accounts are not continued in next months.

For other active accounts in the next month there are generated events: payment or missing payment. It is based on two scorings:

$$\begin{aligned} Score_{Main} &= \sum_{\alpha} \beta_{\alpha}^a x_{\alpha}^a + \sum_{\gamma} \beta_{\gamma}^l x_{\gamma}^l + \sum_{\delta} \beta_{\delta}^{act} x_{\delta}^{act} \\ &+ \sum_{\eta} \sum_t \beta_{\eta}^{beh}(t) x_{\eta}^{beh}(t) + \beta_r \epsilon + \beta_0, \end{aligned} \quad (2.2)$$

$$\begin{aligned} Score_{Cycle} &= \sum_{\alpha} \phi_{\alpha}^a x_{\alpha}^a + \sum_{\gamma} \phi_{\gamma}^l x_{\gamma}^l + \sum_{\delta} \phi_{\delta}^{act} x_{\delta}^{act} \\ &+ \sum_{\eta} \sum_t \phi_{\eta}^{beh}(t) x_{\eta}^{beh}(t) + \phi_r \epsilon + \phi_0, \end{aligned} \quad (2.3)$$

where  $t = 3, 6, 9, 12$ ,  $\alpha = Income, Spending, Nom_1, \dots, Nom_4, Int_1, \dots, Int_4$ ,  $\gamma = Inst, N_{Inst}, Amount$ ,  $\eta = days, n_{due}$ ,  $\delta = days, n_{paid}, n_{due}, utl, dueutl$ ,  $age, capacity, dueinc, loaninc, seniority$ ,  $\varepsilon$  and  $\epsilon$  are taken from the standardized normal distribution  $N$ .

Let consider the following migration matrix:

$$M_{ij}^{act} = \begin{cases} M_{ij}^{adj} & \text{when } Score_{Cycle} \leq \text{Cutoff}, \\ M_{ij} & \text{when } Score_{Cycle} > \text{Cutoff}, \end{cases}$$

where Cutoff is another parameter like all  $\beta$ s and  $\phi$ s.

For fixed  $T_{cur}$  and fixed  $x_{n_{due}}^{act} = i$  all active accounts can be segmented by  $Score_{Main}$  to satisfy the same proportions like appropriate elements of migration matrix  $M_{ij}^{act}$ : the first group  $g = 0$  by the highest scores has share equaled to  $M_{i0}^{act}$ , the second  $g = 1$  has share  $M_{i1}^{act}$ , ..., and the last group  $g = 7$  share –  $M_{i7}^{act}$ .

For particular account assigned to the group  $g$  payment is done in month  $T_{cur} + 1$  when  $g \leq i$ , in other case payment is missing.

For missing payment Transaction dataset is updated by the following information:

$$\begin{aligned} x_{n_{paid}}^t &= x_{n_{paid}}^{act}, \\ x_{n_{due}}^t &= g, \\ x_{days}^t &= \text{Missing}. \end{aligned}$$

For payment by formulas:

$$\begin{aligned} x_{n_{paid}}^t &= \min(x_{n_{paid}}^{act} + x_{n_{due}}^{act} - g + 1, x_{N_{inst}}^l), \\ x_{n_{due}}^t &= g, \end{aligned}$$

and  $x_{days}^t$  are generated from the distribution  $D_{days}$ .

Described steps are repeated for all months between  $T_s$  and  $T_e$ .

## 2.8 Default definition

The Default is a typical credit scoring and Basel II notion. Every account from the observation point  $T_{cur}$  is tested during the outcome period equals 3, 6, 9 and 12 months. During that time there is analyzed maximal number of due installments, exactly:

$$\text{MAX} = \text{MAX}_{m=0}^{t-1}(x_{n_{due}}^{act}(T_{cur} + m)),$$

where  $t = 3, 6, 9, 12$ . Dependently on value MAX are defined three values of default statuses  $\text{Default}_t$ :

**Good:** When  $\text{MAX} \leq 1$  or during the outcome period was  $x_{status}^t = C$ .

**Bad:** When  $\text{MAX} > 3$  or during the outcome period  $x_{status}^t = B$ . In the case  $t = 3$  when  $\text{MAX} > 2$ .

**Indeterminate:** for other cases.

Existing of Indeterminate status can be questionable. In some analysis only two statuses are preferable, for example in Basel II. It is also a good topic for further research which can be solved due to data generator described in this paper.

## 2.9 Portfolio segmentation and risk measures

Typically credit scoring is used for the control of the following sub-portfolios or processes:

**Acceptance process – APP portfolio:** It is the set of all starting points of credits, where it is decided which one are accepted or rejected. Acceptance sub-portfolio is defined as the set of rows of Transaction dataset with the condition:  $T_{cur} = T_{app}$ . Every account belongs to that set only ones.

**Cross-up sell process – BEH portfolio:** It is the set of all accounts with the longer history than 2 months and in the good condition (without delinquency). Cross-up sell or Behavioral sub-portfolio is defined as the set of rows of Transaction dataset with the condition:  $x_{seniority}^{act} > 2$  and  $x_{ndue}^{act} = 0$ . Every account can belongs to that set many times.

**Collection process – COL portfolio:** It is the set of all accounts with the delinquency, but at the beginning of the collection process. Collection sub-portfolio is defined as the set of rows of Transaction dataset with the condition:  $x_{ndue}^{act} = 1$ . Every account can belongs to that set many times.

For every mentioned sub-portfolio one can calculates and tests risk measures called bad rates defined as the share of **Bad** statuses for every observation points and outcome periods.

Definitions of mentioned sub-portfolios in the reality can be more complex, here are suggested the simplest versions for further analysis of cases studies presented in the section 4.

### 3 General theory

#### 3.1 The main assumption and definition

**Definition.** The layout

$$(T_s, T_e, M_{ij}, E(m), \beta_\alpha^a, \beta_\gamma^l, \beta_\delta^{act}, \beta_\eta^{beh}(t), \beta_r, \beta_0,$$

$$\phi_\alpha^a, \phi_\gamma^l, \phi_\delta^{act}, \phi_\eta^{beh}(t), \phi_r, \phi_0, \varepsilon, \epsilon, D_{Birth}, D_\alpha, D_\gamma, D_{Applications}, D_{days}, \text{Cutoff})$$

with the all rules and symbols, relations and processes described in the section 2 is called **The Retail Consumer Finance Data Generator in the case of fixed installment loans** with the nick name **RCFDG**.

**Theorem – assumption.** Every consumer finance portfolio with the fixed installment loans can be estimated by the **RCFDG**.

The proof of that theorem can be always done in the correct way due to parts:  $\beta_r \varepsilon$  and  $\phi_r \epsilon$  in the formulas 2.2 and 2.3. From the empirical point of view credit scoring is always used in portfolio control, so mentioned theorem is correct, but problem is with the goodness of fit. Up to now theory is too early to define a good measures of fit, however it is a proper starting point in the next development of the general theory of consumer finance portfolios.

The similar ideas and researches are presented in [3].

#### 3.2 Open questions

The next steps probably would be concentrated on:

- Finding the correct goodness of fit statistics measuring the distance between the real consumer finance portfolio and **RCFDG**. Also it should be tested the property of that statistics.
- Analyzing the additional constraints to satisfy for example properties like: the predictive power, measured for example by Gini [12], of characteristic  $x_{days}^{beh}(3)$  on  $\text{Default}_6$  should be equalled to 40%.
- Creating more general case with all collection processes, more than one credit per customer, more than one macro-economic factors and other detailed issues.

- Analyzing of various existing real consumer finance portfolios and finding the set of parameters describing each of them. Then there can be developed the theory of principal component analysis (PCA) of all consumer finance portfolios in the particular country or in the world.
- Defining the generalization of the notion of consumer finance portfolio contains almost all properties of real portfolios.
- Using that generalized notion in researches on the development of scoring methods to use that notion as a general idea of method proving. For example the theorem: *Scoring models build on  $Default_3$  and on  $Default_{12}$  produce the same results* could be solved by the additional condition: betas for  $t = 3$  and for  $t = 12$  should be similar. It is very probable that many future researches will discover many properties and relations among betas, coefficients of the migration matrix and their consequences.

## 4 Two case studies

### 4.1 Common parameters

All random numbers are based on two typical random generators: uniform  $U$  and standardized normal  $N$  distributions, in details: the distribution  $U$  returns a number from the interval  $(0, 1)$  with the equal probability.

All common coefficients are the following:  $T_s = 1970.01$  (January 1970),  $T_e = 1976.12$  (December 1976),

$$M_{ij} = \begin{bmatrix} & j=0 & j=1 & j=2 & j=3 & j=4 & j=5 & j=6 & j=7 \\ i=0 & 0.850 & 0.150 & 0.000 & 0.000 & 0.000 & 0.000 & 0.000 & 0.000 \\ i=1 & 0.250 & 0.450 & 0.300 & 0.000 & 0.000 & 0.000 & 0.000 & 0.000 \\ i=2 & 0.040 & 0.240 & 0.190 & 0.530 & 0.000 & 0.000 & 0.000 & 0.000 \\ i=3 & 0.005 & 0.025 & 0.080 & 0.100 & 0.790 & 0.000 & 0.000 & 0.000 \\ i=4 & 0.000 & 0.000 & 0.010 & 0.080 & 0.090 & 0.820 & 0.000 & 0.000 \\ i=5 & 0.000 & 0.000 & 0.000 & 0.000 & 0.020 & 0.030 & 0.950 & 0.000 \\ i=6 & 0.000 & 0.000 & 0.000 & 0.000 & 0.000 & 0.010 & 0.010 & 0.980 \end{bmatrix},$$

$E(m) = 0.01 + (1.5 + \sin((5 \cdot \pi \cdot m)/(T_e - T_s)) + N/5)/8$ ,  $D_{Applications} = 300 \cdot 30 \cdot (1 + N/20)$ , if  $T_{app}$  is December then  $D_{Applications} = D_{Applications} \cdot 1.2$ . To define  $D_{Birth}$  first is defined distribution of age:  $D_{Age} = ((75 - 18) \cdot (N + 4)/7 + 10 + 20 \cdot U)$  if  $Age > 75$  then  $Age = 75$ , if  $Age < 18$  then  $Age = 18$ .  $D_{Birth} = T_{app} - D_{Age} \cdot 365.5$ ,  $D_{Income} = \text{int}((10000 - 500)/40 \cdot 10 \cdot \text{abs}(N) + 500)$ ,  $D_{Inst} = \text{int}(Income \cdot \text{abs}(N)/4)$ ,  $D_{Spending} = \text{int}(Income \cdot \text{abs}(N)/4)$ ,  $D_{N_{Inst}} = \text{int}(30 \cdot \text{abs}(N)/4 + 6)$  if  $N_{Inst} < 6$  then  $N_{Inst} = 6$ ,  $D_{Nom_i} = \text{int}(5 \cdot \text{abs}(N))$  and  $D_{Int_i} = 10 \cdot U$ , for  $i = 1, 2, 3, 4$ , if  $x_{n_{due}}^{act} < 2$  then  $D_{days} = -\text{int}(15 \cdot (\text{abs}(N)/4))$  else  $D_{days} = \text{int}(15 \cdot (N/4))$ , where  $\text{int}()$  and  $\text{abs}()$  are integer value and absolute value suitable.

To avoid scale or unit problem for every individual variable it is suggested to make a simple standardization step for ABT table for every  $T_{cur}$  before score calculation. That idea is quite realistic, because even some customers are good payers in the crisis time they can also have more problems, so general condition of the current month can influence on all customers. On the other hand to present interesting two cases is decided to standardize variables by the global parameters.

Scoring formula for  $Score_{Main}$  is calculated based on the table 1, namely:

$$Score_{Main} = \sum_{index=1}^{28} \beta(x - \mu)/\sigma.$$

All beta coefficients could be recalculated without standardization step, but in that case it would be more difficult to interpret them. By a simple study of the table 1 it can be indicated that the most significant variables have absolute value equals 6.

## 4.2 The first case study – unstable application characteristic – APP

In that case it is assumed that only customers with low income can be influenced by a crisis. Application characteristic income in that data generator is a stable variable during the time, and the migration matrix is adjusted by the macro-economic  $E(m)$  only for cases:

$$x_{Income}^a < 1800.$$

Presented relation without any problem can be transformed into the general form 2.3.

Table 1: Scoring formula for  $Score_{Main}$ .

Index	$x$ – variable	$\mu$	$\sigma$	$\beta$
1	$x_{Nom_1}^a$	3.5	3	1
2	$x_{Nom_2}^a$	3.5	3	2
3	$x_{Nom_3}^a$	3.5	3	1
4	$x_{Nom_4}^a$	3.5	3	3
5	$x_{Int_1}^a$	5	2.89	1
6	$x_{Int_2}^a$	5	2.89	-4
7	$x_{Int_3}^a$	5	2.89	1
8	$x_{Int_4}^a$	5	2.89	-2
9	$x_{days}^{act}$	13	2.42	-5
10	$x_{utl}^{act}$	0.36	0.28	-4
11	$x_{du\epsilon utl}^{act}$	0.12	0.2	-6
12	$x_{n_{due}}^{act}$	1.3	2	-2
13	$x_{age}^{act}$	53	9.9	4
14	$x_{capacity}^{act}$	0.4	0.21	-2
15	$x_{dueinc}^{act}$	0.3	0.6	-1
16	$x_{loaninc}^{act}$	2.4	2.1	-2
17	$x_{income}^a$	2395	1431	2
18	$x_{Amount}^a$	5741	6804	-1
19	$x_{N_{inst}}^l$	12.3	4.63	-4
20	$x_{due}^{beh}(3)$	1.4	1.6	-4
21	$x_{days}^{beh}(3)$	14.15	1.4	-6
22	$x_{due}^{beh}(6)$	1.6	1.13	-5
23	$x_{days}^{beh}(6)$	14.57	1.02	-6
24	$x_{due}^{beh}(9)$	1.78	0.75	-5
25	$x_{days}^{beh}(9)$	14.78	0.72	-6
26	$x_{due}^{beh}(12)$	1.89	0.48	-5
27	$x_{days}^{beh}(12)$	14.91	0.49	-6
28	$\epsilon$	0	0.02916	1

### 4.3 The second case study – unstable behavioral characteristic – BEH

Here the condition for migration matrix adjustment is the following:

$$x_{n_{due}}^{beh}(6) > 0 \quad \text{and} \quad x_{seniority}^{act} > 6,$$

the rule for the seniority variable is added to not adjust accounts with missing imputation based on 2.1. That case presents situation when crisis has an impact on customers who had some delinquency during their last 6 months.

### 4.4 Stability problem

Let be considered the typical scoring models building process, for example on behavioral sub-portfolio. Because two cases are based on two variables one application and one behavioral let be considered only the set of these two variables. To indicate strong instability models they are analyzed with the target variable Default<sub>9</sub>.

Every variable is segmented or binned for a few attributes described in the tables 2 and 3.

In the case of unstable application variable (APP) by studying the figure 6 can be confirmed, what is expected, that attribute 2 is very stable during the time and accounts from that group are not quite sensible for crisis changes. In opposite attribute 1 is very unstable. The same groups in the case of unstable behavioral variable (BEH) are both unstable, see the figure 7. The same group, accounts from attribute 2, are presented on figure 5 for both cases to indicate in a better scale that APP case can really choose accounts not sensitive on the crisis. Even data generator is simplicity of the real data, that conclusion is very useful. Some application data can be profitable in risk management to indicate sub-segments with stable risk in the time.

Not the same conclusions can be formulated for behavioral variable  $x_{n_{due}}^{beh}(6)$ . On the figure 3 there are presented risk evolutions for three attributes of that variable. All of them are not stable. The most stable attribute is with the number 3. Also for the case BEH that attribute is not stable, see the figure 4. To be sure of that there are also presented on the figure 2 only attributes 3 for both cases. Every reader can say that both cases have unstable risk. Even in the case BEH the attribute 3 is expected to have a stable risk, due to the rule for migration matrix adjustment, expectation has failed. The reason comes from the correct understanding of the process. Typical scoring



Table 2: Simple binning for two variables in the case APP.

Characteristic	Attribute number	Condition	Bad rate on Default <sub>9</sub>	Population percent	Gini on Default <sub>9</sub>
$x_{n_{due}}^{beh}(6)$	1	$x_{seniority}^{act} < 6$	16.77%	37.09%	51.34%
	2	$x_{n_{due}}^{beh}(6) > 0$ and $x_{seniority}^{act} \geq 6$	6.48%	22.49%	
	3	otherwise	1.07%	40.42%	
$x_{Income}^a$	1	$x_{Income}^a < 1800$	20.11%	18.32%	36.29%
	2	$x_{Income}^a \geq 1800$	4.72%	81.68%	

approach is based on the principal idea that historical information up to the observation point is able to predict behavior during the outcome period. Up to the observation point account did not have any delinquency so the variable  $x_{n_{due}}^{beh}(6) = 0$ . After that point in the next months account can have due installments. It can be adjusted by the macro-economic variable and on the end that group can become unstable.

The mentioned idea is very important for father research of the crisis. It should be emphasized that typical scoring methods used on three types of sub-portfolios: APP, BEH and COL cannot discover in the correct way the rule of crisis adjustment and cannot indicate some sub-segments stable in the time. Of course scoring can be also used just like in that paper for prediction of migration states; to be very clear, not for default statuses prediction but for transition prediction. The best method is probably the survival analysis [13] or [14] with time covariates (time dependent variables), where in natural way there is indicated the factor of being better or worse payer in the correct time, namely in the typical scoring model the factor is considered but only up to the observation point. In the survival model however it can be also taken into the account after that observation point, so in the more realistic way.

There are made many other cases of data generators with more complex rule for  $Score_{Cycle}$ . If there are taken together both types of variables: application and behavioral the case is too complicated and unstable property exists everywhere. In that case is not possible to find stable factor. That conclusion is also very important for crisis analysis, because it describes the nature of crisis: if it is a strong event and it has an impact on both types of characteristics behavioral and application – it is and risk management can try to find some sub-segments only more stable then others or with maximal risk not exceeded the expected boundary.

Table 3: Simple binning for two variables in the case BEH.

Characteristic	Attribute number	Condition	Bad rate on Default <sub>9</sub>	Population percent	Gini on Default <sub>9</sub>
$x_{n_{due}}^{beh}(6)$	1	$x_{seniority}^{act} < 6$	19.49%	40.05%	46.54%
	2	$x_{n_{due}}^{beh}(6) > 0$ and $x_{seniority}^{act} \geq 6$	14.04%	16.52%	
	3	otherwise	1.74%	43.43%	
$x_{Income}^a$	1	$x_{Income}^a < 1800$	12.09%	39.49%	5.04%
	2	$x_{Income}^a \geq 1800$	10.09%	60.51%	

## 4.5 Various types of risk measures

Let be defined that crisis is a time where risk is the highest. The most popular reporting for risk management is based on bad rates, vintage and flow rates. The figure 1 presents bad rates for three different sub-portfolios application, behavioral and collection. There is presented also one flow rate. There is a simple conclusion that crisis does not occur in the same time. Some curves indicate local maximum of risk earlier than others. The difference in the time is significant and can be almost 6 months, so it is very important to remember what kind of reports can indicate a crisis as quickly as possible. It should be emphasized that bad rates reports present, by the standard way, the evaluation of risk by observation points and a crisis time can occur between observation point and the end of outcome period. It seems that flow rates reports precise the crisis time in better way.

## 4.6 Implementation

All data were prepared by the SAS System [11] by manual codes written in SAS 4GL used units: Base SAS and SAS/STAT. For the case of unstable behavioral variable – BEH: Production dataset has 779 993 rows (about 90MB) and Transaction dataset – 8 969 413 rows (about 400MB). Total time of calculation per one case takes about 4 hours.

## 5 Conclusions

Even if data are generated by random-simulated process, which is not realistic, the conclusions give the possibility to better understand the nature

Figure 1: Risk measures on  $\text{Default}_9$  comparison on sub-portfolios: APP, BEH and COL and also with one flow rate  $M_{23}$ .

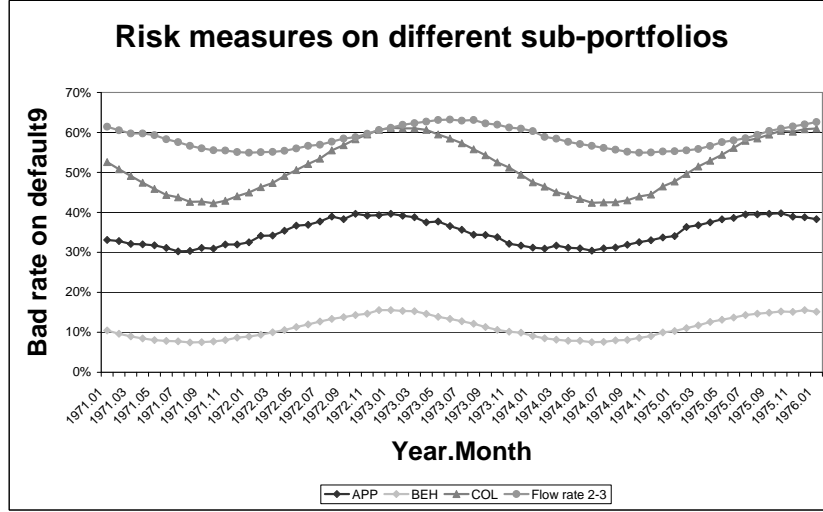


Figure 2: Risk measures on  $\text{Default}_9$  on attribute 3 of variable  $x_{ndue}^{beh}(6)$  for two cases APP and BEH.

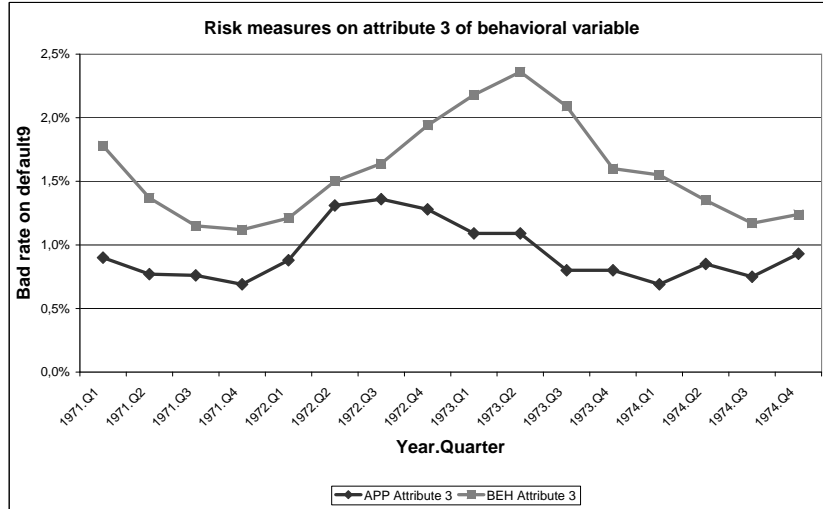


Figure 3: Risk measures on Default<sub>9</sub> on attributes of variable  $x_{n_{due}}^{beh}$  (6) for the case APP.

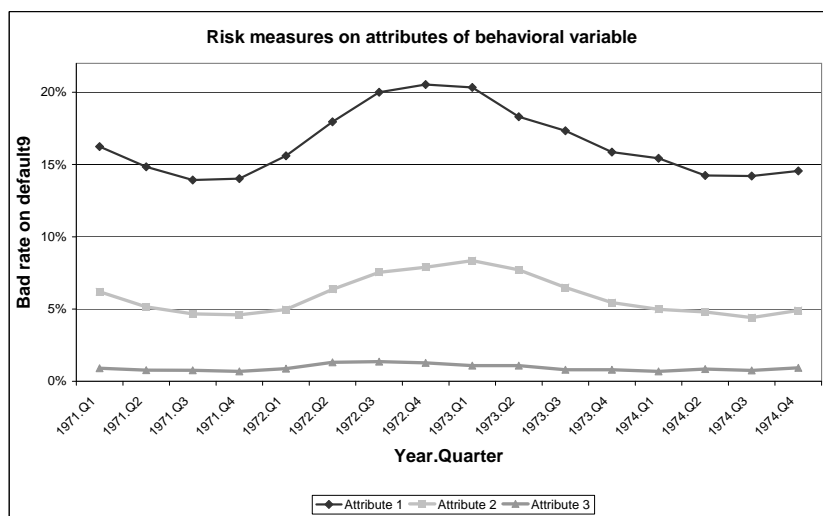


Figure 4: Risk measures on Default<sub>9</sub> on attributes of variable  $x_{n_{due}}^{beh}$  (6) for the case BEH.

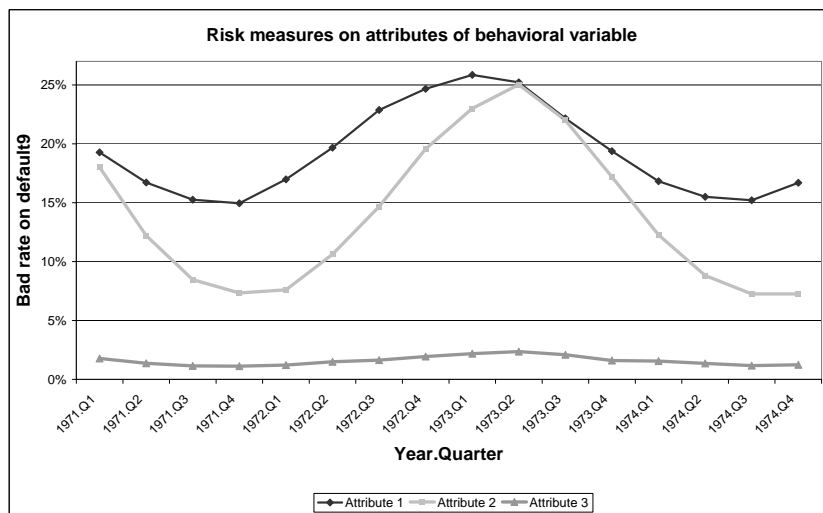


Figure 5: Risk measures on  $\text{Default}_9$  on attribute 2 of variable  $x_{Income}^a$  for two cases APP and BEH.

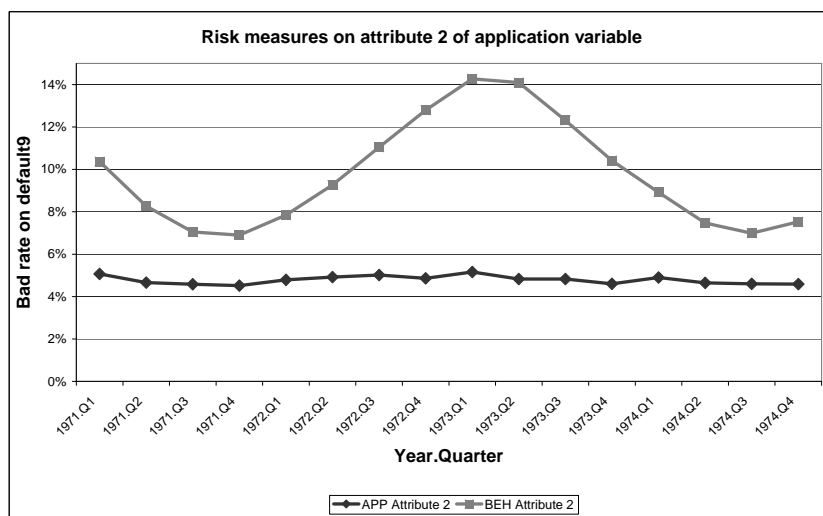


Figure 6: Risk measures on  $\text{Default}_9$  on attributes of variable  $x_{Income}^a$  for the case APP.

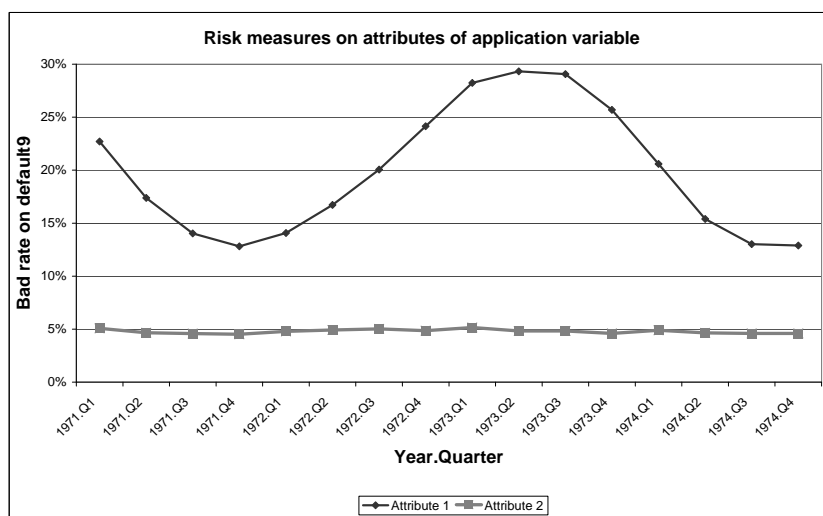
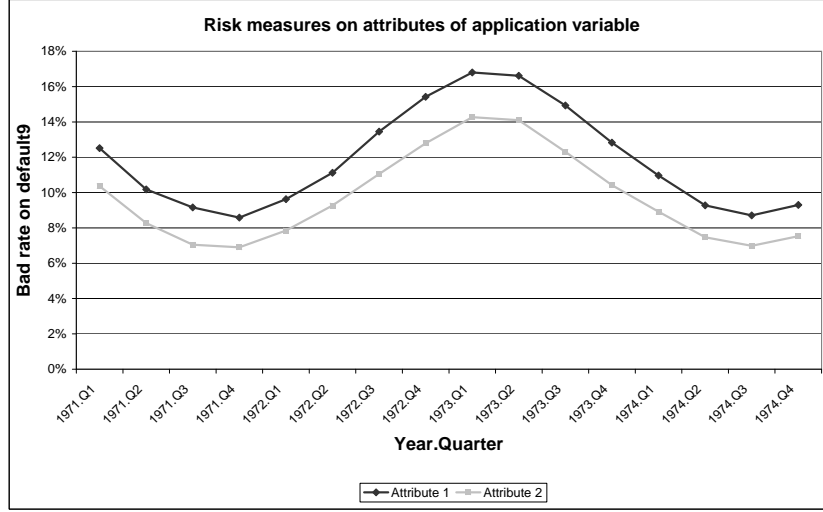


Figure 7: Risk measures on Default<sub>9</sub> on attributes of variable  $x_{Income}^a$  for the case BEH.



of the crisis.

The banking data generator is a new hope for researching to find the proving method of comparisons of various credit scoring techniques. It is probable that in the future many random generated data will become the new repository for testing and comparisons.

In the first case – unstable application variable like income is possible to split portfolio for two parts: stable and unstable during the time. For the second case unstable – behavioral characteristic the task is more complicated and it is not possible to split in the same way. Some sub-segments can have better stability but always they fluctuate. Moreover if a crisis is impacted by many factors both from application form customer characteristics and from a customer behavioral together it is very difficult to indicate these factors and the crisis in reports is everywhere.

Generated data are very useful for various analysis and researches. There are many rows, many bad default statuses, so analyst can make many good exercises to improve his experience.

## References

- [1] Edward Huang. Scorecard specification, validation and user acceptance: A lesson for modellers and risk managers. *Credit Scoring Conference CRC, Edinburgh*, 2007.
- [2] Basel Committee on Banking Supervision. International convergence of capital measurement and capital standards. *A Revised Framework*, Updated November 2005. <http://www.bis.org>.
- [3] Madhur Malik & Lyn C Thomas. Modelling credit risk in portfolios of consumer loans: Transition matrix model for consumer credit ratings. *Credit Scoring Conference CRC, Edinburgh*, 2009.
- [4] Edward Huang, Christopher Scott. Credit risk scorecard design, validation and user acceptance: A lesson for modellers and risk managers. *Credit Scoring Conference CRC, Edinburgh*, 2007.
- [5] Izabela Majer. Application scoring: logit model approach and the divergence method compared. *Warsaw School of Economics – SGH, Working Paper No. 10-06*, 2010.
- [6] Elizabeth Mays. Systematic risk effects on consumer lending products. *Credit Scoring Conference CRC, Edinburgh*, 2009.
- [7] Naeem Siddiqi. Credit risk scorecards: Developing and implementing intelligent credit scoring. *Wiley and SAS Business Series*, 2005.
- [8] OPNET Technologies Inc. <http://www.opnet.com>.
- [9] Bala Supramaniam, Mahadevan; Shanmugam. Simulating retail banking for banking students. *Reports – Evaluative, Practitioners and Researchers ERIC Identifier: ED503907*, 2009.
- [10] H. J. Watson. Simulating retail banking for banking students. *Computer simulation in business. New York: John Wiley & Sons.*, 1981.
- [11] SAS Institute Inc. <http://www.sas.com>.
- [12] Basel Committee on Banking Supervision. Validation of internal rating systems. *Working Paper No. 14*, February 2005. <http://www.bis.org>.

- [13] Tony Bellotti, Jonathan Crook. Credit scoring with macroeconomic variables using survival analysis. *Journal of the Operational Research Society* Key: citeulike:4083586, 2009.
- [14] Jonathan Crook. Dynamic consumer risk models: an overview. *Credit Scoring Conference CRC, Edinburgh*, 2008.